
Lecture6(Part3)

Topics covered:
Memory subsystem



Virtual memories

- ❑ Recall that an important challenge in the design of a computer system is to provide a large, fast memory system at an affordable cost.
- ❑ Architectural solutions to increase the effective speed and size of the memory system.
- ❑ Cache memories were developed to increase the effective speed of the memory system.
- ❑ Virtual memory is an architectural solution to increase the effective size of the memory system.



Virtual memories (contd..)

- ❑ Recall that the addressable memory space depends on the number of address bits in a computer.
 - ◆ For example, if a computer issues 32-bit addresses, the addressable memory space is 4G bytes (divided into user space and system space).
- ❑ Physical main memory in a computer is generally not as large as the entire possible addressable space.
 - ◆ Physical memory typically ranges from a few hundred megabytes to 1G bytes.
- ❑ Large programs that cannot fit completely into the main memory have their parts stored on secondary storage devices such as magnetic disks.
 - ◆ Pieces of programs must be transferred to the main memory from secondary storage before they can be executed.

If you have 2 GB RAM and 1GB virtual memory. A user program needs only 2GB to run, the program cannot run, why????



Virtual memories (contd..)

- ❑ When a new piece of a program is to be transferred to the main memory, and the main memory is full, then some other piece in the main memory must be replaced.
 - ◆ Recall this is very similar to what we studied in case of cache memories.
- ❑ Operating system automatically transfers data between the main memory and secondary storage.
 - ◆ Application programmer need not be concerned with this transfer.
 - ◆ Also, application programmer does not need to be aware of the limitations imposed by the available physical memory.

◆ Virtual memories (contd..)

- ❑ Techniques that automatically move program and data between main memory and secondary storage when they are required for execution are called virtual-memory techniques.
- ❑ Programs and processors reference an instruction or data independent of the size of the main memory.
- ❑ Processor issues binary addresses for instructions and data.
 - ◆ These binary addresses are called logical or virtual addresses.
- ❑ Virtual addresses are translated into physical addresses by a combination of hardware and software subsystems.
 - ◆ If virtual address refers to a part of the program that is currently in the main memory, it is accessed immediately.
 - ◆ If the address refers to a part of the program that is not currently in the main memory, it is first transferred to the main memory before it can be used.

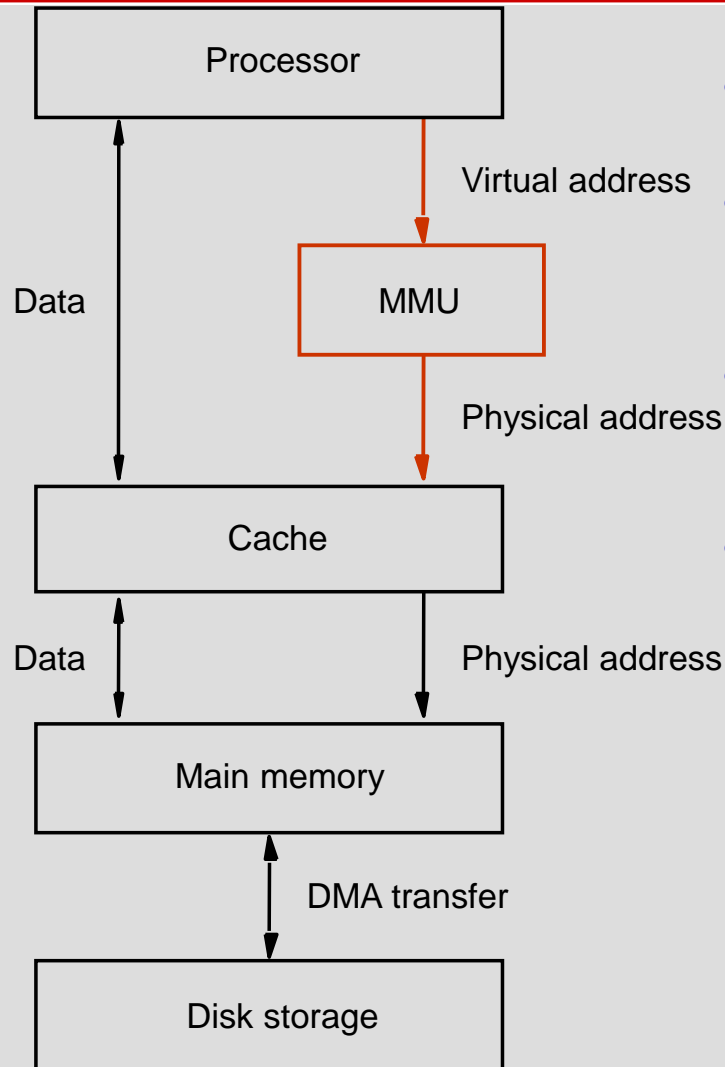
Virtual memory

- ❑ **Virtual memory** is a feature of an operating system (OS) (implemented using software and hardware) that allows a computer to compensate for shortages of physical **memory** by temporarily transferring pages of data from random access **memory** (RAM) to disk storage.
- ❑ In effect, RAM acts like cache for disk.
- ❑ The primary benefits of virtual memory include freeing applications from having to manage a shared memory space, increased security (memory protection) due to memory isolation, and being able to conceptually use more memory than might be physically available, using the technique of paging.
- ❑ <http://www.cs.umd.edu/class/sum2003/cmsc311/Notes/Memory/virtual.html>

◆ Virtual(Logical) address

- ❑ A virtual address is a binary number in virtual memory that enables a process to use a location in primary storage (RAM) independently of other processes and to use more space than actually exists in primary storage (RAM) by temporarily relegating (take away) some contents to a hard disk or internal flash drive.
- ❑ In a computer with both physical and virtual memory, a so-called MMU (memory management unit) coordinates and controls all of the memory resources, assigning portions called pages (large blocks) to various running programs to optimize system performance. By translating between virtual addresses and physical addresses, the MMU allows every running process to "think" that it has all the primary storage (RAM) to itself.

Virtual memory organization



- *Memory management unit (MMU) translates virtual addresses into physical addresses.*
- *If the desired data or instructions are in the main memory they are fetched as described previously.*
- *If the desired data or instructions are not in the main memory, they must be transferred from secondary storage to the main memory.*
- *MMU causes the operating system to bring the data from the secondary storage into the main memory.*



Address translation

- ❑ Assume that program and data are composed of fixed-length units called pages.
- ❑ A page consists of a large block of words that occupy contiguous (neighboring) locations in the main memory.
- ❑ Page is a basic unit of information that is transferred between secondary storage and main memory.
- ❑ Size of a page commonly ranges from 2K to 16K bytes.
 - ◆ Pages should **not be too small**, because data can be transferred at high rates (megabytes per second) between a secondary storage device and the main memory.
 - ◆ Pages should **not be too large**, else a large portion of the page may not be used, and it will occupy valuable space in the main memory.



Address translation (contd..)

- ❑ Concepts of virtual memory are similar to the concepts of cache memory.
- ❑ Cache memory:
 - ◆ Introduced to bridge the speed gap between the processor and the main memory.
 - ◆ Implemented in hardware.
- ❑ Virtual memory:
 - ◆ Introduced to bridge the speed gap between the main memory and secondary storage.
 - ◆ Implemented in hardware and software.



Address translation (contd..)

- ❑ Each virtual or logical address generated by a processor is interpreted as a virtual page number (high-order bits) plus an offset (low-order bits) that specifies the location of a particular byte within that page.
- ❑ Information about the main memory location of each page is kept in the page table.
 - ◆ Main memory address where the page is stored.
 - ◆ Current status of the page.
- ❑ Area of the main memory that can hold a page is called as page frame.
- ❑ Starting address of the page table is kept in a page table base register.

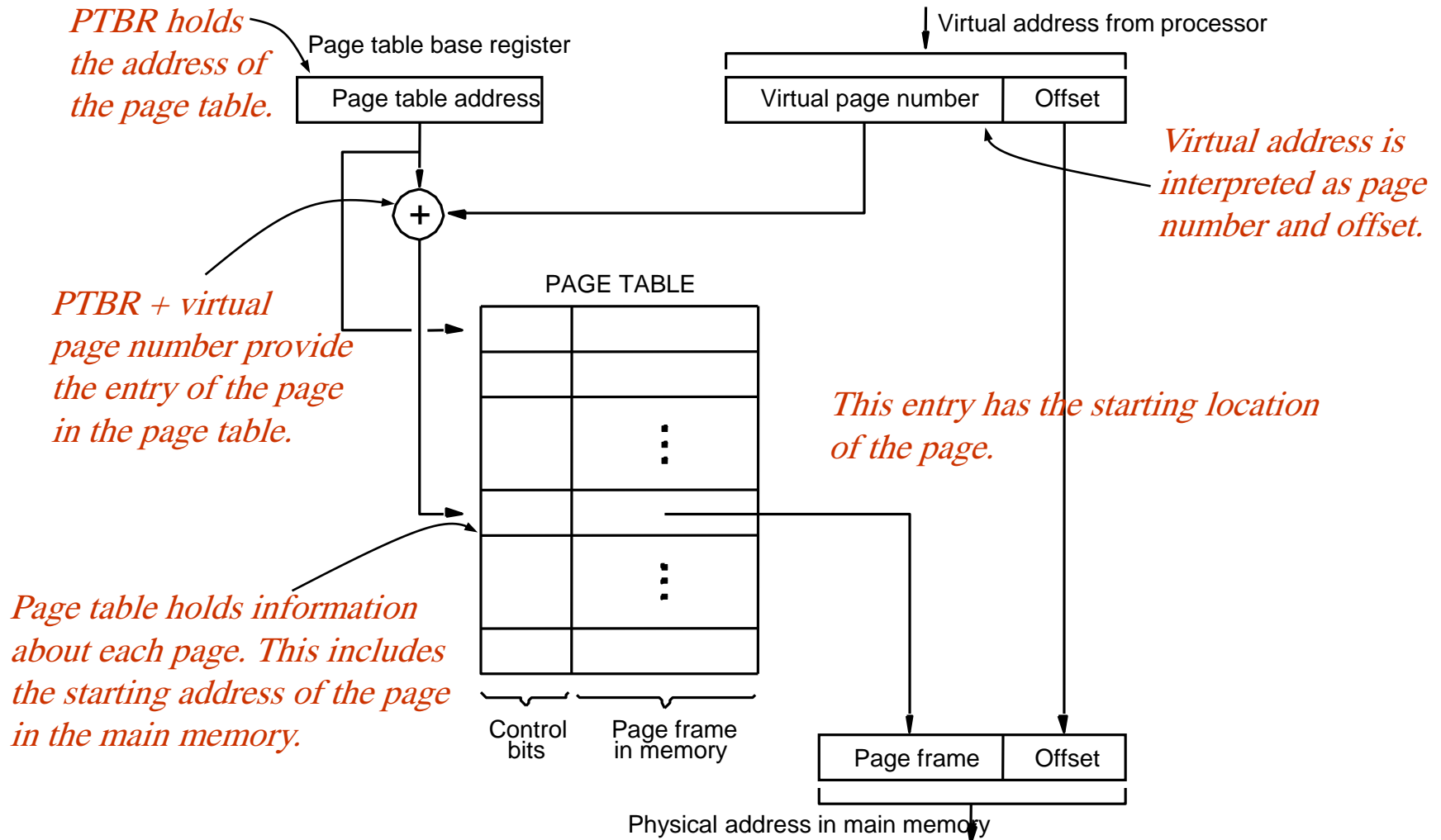


Address translation (contd..)

- ❑ Virtual page number generated by the processor is added to the contents of the page table base register.
 - ◆ This provides the address of the corresponding entry in the page table.
- ❑ The contents of this location in the page table give the starting address of the page if the page is currently in the main memory.



Address translation (contd..)





Address translation (contd..)

- ❑ Page table entry for a page also includes some control bits which describe the status of the page while it is in the main memory.
- ❑ One bit indicates the validity of the page.
 - ◆ Indicates whether the page is actually loaded into the main memory.
 - ◆ Allows the operating system to invalidate the page without actually removing it.
- ❑ One bit indicates whether the page has been modified during its residency in the main memory.
 - ◆ This bit determines whether the page should be written back to the disk when it is removed from the main memory.
 - ◆ Similar to the dirty or modified bit in case of cache memory.



Address translation (contd..)

- Other control bits for various other types of restrictions that may be imposed.
 - ◆ For example, a program may only have read permission for a page, but not write or modify permissions.

For a 4K page you require $(4K == (4 * 1024) == 4096 == 2^{12} ==)$ 12 bits of offset.

32-bit address-space would require a table of 1048576 entries when using 4KB pages.



Address translation (contd..)

- ❑ Where should the page table be located?
- ❑ Recall that the page table is used by the MMU for every read and write access to the memory.
 - ◆ Ideal location for the page table is within the MMU.
- ❑ Page table is quite large.
- ❑ MMU is implemented as part of the processor chip.
- ❑ Impossible to include a complete page table on the chip.
- ❑ Page table is kept in the main memory.
- ❑ A copy of a small portion of the page table can be accommodated within the MMU.
 - ◆ Portion consists of page table entries that correspond to the most recently accessed pages.



End